**Amendment to the Claims**

This listing of claims will replace all prior versions, and listings, of claims in the application:

**Listing of Claims:**

1       1.       (currently amended) A system for grouping clusters of
2    semantically scored documents electronically stored in a data corpus, comprising:
3                a scoring module determining a score, which is assigned to at least one
4    concept that has been extracted from a plurality of electronically-stored
5    documents, wherein the score is calculated as a function of a summation of a
6    frequency of occurrence of the at least one concept within at least one such
7    document, a concept weight based on a number of terms for the at least one
8    concept, a structural weight, and a corpus weight, forming the score assigned to
9    the at least one concept as a normalized score vector for each such document, and
10   determining a similarity between the normalized score vector for each such
11   document as an inner product of each normalized score vector;
12               a clustering module forming clusters of the documents, comprising:
13                       a selection submodule selecting a set of candidate seed documents
14   from the plurality of documents;
15                       a seed document identification submodule identifying a set of seed
16   documents by applying the similarity to each such candidate seed document and
17   selecting those candidate seed documents that are sufficiently unique from other
18   candidate seed documents as the seed documents;
19                       a non-seed document identification submodule identifying a
20   plurality of non-seed documents;
21                       a comparison submodule determining the similarity between each
22   non-seed document and a center of each cluster; and
23                       a clustering submodule grouping each such non-seed document
24   into a cluster with a best fit, subject to a minimum fit;

25        a threshold module determining the similarity between each of the

26    documents grouped into each cluster based on the center of the cluster and the

27    scores assigned to each of the at least one concepts in that document, dynamically

28    determining a threshold for each cluster as a function of the similarity between

29    each of the documents, and identifying and reassigning each of the documents

30    having the similarity falling outside the threshold.

1        Claim 2 (canceled).

1        3.      (previously presented) A system according to Claim 1, further

2    comprising:

3        a compression module compressing the score through logarithmic

4    compression.

1        Claim 4 (canceled).

1        5.      (original) A system according to Claim 1, further comprising:

2        the scoring module calculating the structural weight as a function of a

3    location of the at least one concept within the at least one such document.

1        6.      (original) A system according to Claim 1, further comprising:

2        the scoring module calculating the corpus weight as a function of a

3    reference count of the at least one concept over the plurality of documents.

1        Claims 7-8 (canceled).

1        9.      (currently amended) A method for grouping clusters of

2    semantically scored documents electronically stored in a data corpus, comprising:

3        determining a score, which is assigned to at least one concept that has

4    been extracted from a plurality of electronically-stored documents, wherein the

5    score is calculated as a function of a summation of a frequency of occurrence of

6    the at least one concept within at least one such document, a concept weight <u>based</u>

7 on a number of terms for the at least one concept, a structural weight, and a

8 corpus weight;

9    forming the score assigned to the at least one concept as a normalized

10 score vector for each such document;

11    determining a similarity between the normalized score vector for each

12 such document as an inner product of each normalized score vector;

13    forming logically-grouped clusters of the documents, comprising:

14      selecting a set of candidate seed documents from the plurality of

15 documents;

16      identifying a set of seed documents by applying the similarity to

17 each such candidate seed document and selecting those candidate seed documents

18 that are sufficiently unique from other candidate seed documents as the seed

19 documents;

20      identifying a plurality of non-seed documents;

21      determining the similarity between each non-seed document and a

22 center of each cluster; and

23      grouping each such non-seed document into a cluster with a best

24 fit, subject to a minimum fit;

25    determining the similarity between each of the documents grouped into

26 each cluster based on the center of the cluster and the scores assigned to each of

27 the at least one concepts in that document;

28    dynamically determining a threshold for each cluster as a function of the

29 similarity between each of the documents; and

30    identifying and reassigning each of the documents having the similarity

31 falling outside the threshold.

1    Claim 10 (canceled).

1    11.  (previously presented) A method according to Claim 9, further

2 comprising:

3    compressing the score through logarithmic compression.

1        Claim 12 (canceled).

1        13.    (original) A method according to Claim 9, further comprising:

2        calculating the structural weight as a function of a location of the at least

3  one concept within the at least one such document.

1        14.    (original) A method according to Claim 9, further comprising:

2        calculating the corpus weight as a function of a reference count of the at

3  least one concept over the plurality of documents.

1        Claims 15-16 (canceled).

1        17.    (currently amended) A computer-readable storage medium holding

2  code for grouping clusters of semantically scored documents electronically stored

3  in a data corpus, comprising:

4        code for determining a score, which is assigned to at least one concept that

5  has been extracted from a plurality of electronically-stored documents, wherein

6  the score is calculated as a function of a summation of a frequency of occurrence

7  of the at least one concept within at least one such document, a concept weight

8  <u>based on a number of terms for the at least one concept</u>, a structural weight, and a

9  corpus weight;

10       code for forming the score assigned to the at least one concept as a

11  normalized score vector for each such document;

12       code for determining a similarity between the normalized score vector for

13  each such document as an inner product of each normalized score vector;

14       code for forming logically-grouped clusters of the documents, comprising;

15           code for selecting a set of candidate seed documents from the

16  plurality of documents;

17           code for identifying a set of seed documents by applying the

18  similarity to each such candidate seed document and selecting those candidate

19    seed documents that are sufficiently unique from other candidate seed documents

20    as the seed documents;

21            code for identifying a plurality of non-seed documents;

22            code for determining the similarity between each non-seed

23    document and a center of each cluster; and

24            code for grouping each such non-seed document into a cluster with

25    a best fit, subject to a minimum fit;

26            code for determining the similarity between each of the documents

27    grouped into each cluster based on the center of the cluster and the scores

28    assigned to each of the at least one concepts in that document;

29            code for dynamically determining a threshold for each cluster as a

30    function of the similarity between each of the documents; and

31            code for identifying and reassigning each of the documents having the

32    similarity falling outside the threshold.


1        18.    (currently amended) A system for providing efficient document

2    scoring of concepts within and clustering of documents in an electronically-stored

3    document set, comprising:

4            a scoring module scoring a document in an electronically-stored document

5    set, comprising:

6                a frequency module determining a frequency of occurrence of at

7    least one concept within a document;

8                a concept weight module analyzing a concept weight reflecting a

9    specificity of meaning for the at least one concept within the document, wherein

10    the concept weight is based on a number of terms for the at least one concept;

11                a structural weight module analyzing a structural weight reflecting

12    a degree of significance based on structural location within the document for the

13    at least one concept;

14                a corpus weight module analyzing a corpus weight inversely

15    weighing a reference count of occurrences for the at least one concept within the

16    document;

17          a scoring evaluation module evaluating a score to be associated

18  with the at least one concept as a function of a summation of the frequency,

19  concept weight, structural weight, and corpus weight;

20          a vector module forming the score assigned to the at least one

21  concept as a normalized score vector for each such document in the

22  electronically-stored document set; and

23          a determination module determining a similarity between the

24  normalized score vector for each such document as an inner product of each

25  normalized score vector;

26      a clustering module grouping the documents by the score into a plurality

27  of clusters, comprising:

28          a selection submodule selecting a set of candidate seed documents

29  from the electronically-stored document set;

30          a cluster seed submodule identifying seed documents by applying

31  the similarity to each such candidate seed document and selecting those candidate

32  seed documents that are sufficiently unique from other candidate seed documents

33  as the seed documents;

34          an identification submodule identifying a plurality of non-seed

35  documents;

36          a comparison submodule determining the similarity between each

37  non-seed document and a cluster center of each cluster; and

38          a clustering submodule assigning each non-seed document to the

39  cluster with a best fit, subject to a minimum fit; and

40      a threshold module relocating outlier documents, comprising determining

41  the similarity between each of the documents grouped into each cluster based on

42  the center of the cluster and the scores assigned to each of the at least one

43  concepts in that document, dynamically determining a threshold for each cluster

44  as a function of the similarity between each of the documents, and identifying and

45  reassigning each of the documents with the similarity falling outside the

46  threshold.

1    19.    (previously presented) A system according to Claim 18, further

2  comprising:

3       the scoring module evaluating the score in accordance with the formula:

4  $$S_i = \sum_{1 \to n}^{j} f_{ij} \times cw_{ij} \times sw_{ij} \times rw_{ij}$$

5  where $S_i$ comprises the score, $f_{ij}$ comprises the frequency, $0 < cw_{ij} \leq 1$ comprises

6  the concept weight, $0 < sw_{ij} \leq 1$ comprises the structural weight, and $0 < rw_{ij} \leq 1$

7  comprises the corpus weight for occurrence $j$ of concept $i$.

1    20.    (currently amended) A system according to Claim 19, further

2  comprising:

3       the concept weight module evaluating the concept weight in accordance

4  with the formula:

5  $$cw_{ij} = \begin{cases} 0.25 + (0.25 \times t_{ij}), & 1 \leq t_{ij} \leq 3 \\ 0.25 + (0.25 \times [7 - t_{ij}]), & 4 \leq t_{ij} \leq 6 \\ 0.25, & t_{ij} \geq 7 \end{cases}$$

6  where $cw_{ij}$ comprises the concept weight and $t_{ij}$ comprises [[a]] the number of

7  terms for occurrence $j$ of each such concept $i$.

1    21.    (previously presented) A system according to Claim 19, further

2  comprising:

3       the structural weight module evaluating the structural weight in

4  accordance with the formula:

5  $$sw_{ij} = \begin{cases} 1.0, & if(j \approx SUBJECT) \\ 0.8, & if(j \approx HEADING) \\ 0.7, & if(j \approx SUMMARY) \\ 0.5 & if(j \approx BODY) \\ 0.1 & if(j \approx SIGNATURE) \end{cases}$$

6  where $sw_{ij}$ comprises the structural weight for occurrence $j$ of each such concept $i$.

1        22.      (previously presented) A system according to Claim 19, further

2    comprising:

3        the corpus weight module evaluating the corpus weight in accordance with

4    the formula:

5
$$rw_{ij} = \begin{cases} \left(\dfrac{T-r_{ij}}{T}\right)^2, & r_{ij} > M \\ 1.0, & r_{ij} \leq M \end{cases}$$

6    where $rw_{ij}$ comprises the corpus weight, $r_{ij}$ comprises a reference count for

7    occurrence $j$ of each such concept $i$, $T$ comprises a total number of reference

8    counts of documents in the document set, and $M$ comprises a maximum reference

9    count of documents in the document set.

1        23.      (previously presented) A system according to Claim 19, further

2    comprising:

3        a compression module compressing the score in accordance with the

4    formula:

5        $$S_i' = \log(S_i + 1)$$

6    where $S_i'$ comprises the compressed score for each such concept $i$.

1        24.      (original) A system according to Claim 18, further comprising:

2        a global stop concept vector cache maintaining concepts and terms; and

3        a filtering module filtering selection of the at least one concept based on

4    the concepts and terms maintained in the global stop concept vector cache.

1        25.      (original) A system according to Claim 18, further comprising:

2        a parsing module identifying terms within at least one document in the

3    document set, and combining the identified terms into one or more of the

4    concepts.

1        26.      (original) A system according to Claim 25, further comprising:

2      the parsing module structuring each such identified term in the one or

3    more concepts into canonical concepts comprising at least one of word root,

4    character case, and word ordering.

1      27.    (original) A system according to Claim 25, wherein at least one of

2    nouns, proper nouns and adjectives are included as terms.

1      Claims 28-30 (canceled).

1      31.    (previously presented) A system according to Claim 18, further

2    comprising:

3      the similarity submodule calculating the similarity in accordance with the

4    formula:

5
$$\cos \sigma_{AB} = \frac{\langle \vec{S}_A \cdot \vec{S}_B \rangle}{|\vec{S}_A||\vec{S}_B|}$$

6    where $\cos \sigma_{AB}$ comprises a similarity between a document $A$ and a document $B$,

7    $\vec{S}_A$ comprises a score vector for document $A$, and $\vec{S}_B$ comprises a score vector for

8    document $B$.

1      Claims 32-34 (canceled).

1      35.    (currently amended) A method for providing efficient document

2    scoring of concepts within and clustering of documents in an electronically-stored

3    document set, comprising:

4      scoring a document in an electronically-stored document set, comprising:

5          determining a frequency of occurrence of at least one concept

6    within a document;

7          analyzing a concept weight reflecting a specificity of meaning for

8    the at least one concept within the document, wherein the concept weight is based

9    on a number of terms for the at least one concept;

10          analyzing a structural weight reflecting a degree of significance

11  based on structural location within the document for the at least one concept;

12          analyzing a corpus weight inversely weighing a reference count of

13  occurrences for the at least one concept within the document; and

14          evaluating a score to be associated with the at least one concept as

15  a function of a summation of the frequency, concept weight, structural weight,

16  and corpus weight;

17          forming the score assigned to the at least one concept as a normalized

18  score vector for each such document in the electronically-stored document set;

19          determining a similarity between the normalized score vector for each

20  such document as an inner product of each normalized score vector;

21          grouping the documents by the score into a plurality of clusters,

22  comprising:

23          selecting a set of candidate seed documents from the

24  electronically-stored document set;

25          identifying seed documents by applying the similarity to each such

26  candidate seed document and selecting those candidate seed documents that are

27  sufficiently unique from other candidate seed documents as the seed documents;

28          identifying a plurality of non-seed documents;

29          determining the similarity between each non-seed document and a

30  center of each cluster; and

31          assigning each non-seed document to the cluster with a best fit,

32  subject to a minimum fit; and

33          relocating outlier documents, comprising:

34          determining the similarity between each of the documents grouped

35  into each cluster based on the center of the cluster and the scores assigned to each

36  of the at least one concepts in that document;

37          dynamically determining a threshold for each cluster as a function

38  of the similarity between each of the documents; and

39              identifying and reassigning each of the documents with the

40   similarity falling outside the threshold.

1        36.       (previously presented) A method according to Claim 35, further

2   comprising:

3            evaluating the score in accordance with the formula:

4
$$S_i = \sum_{1 \to n}^{j} f_{ij} \times cw_{ij} \times sw_{ij} \times rw_{ij}$$

5   where $S_i$ comprises the score, $f_{ij}$ comprises the frequency, $0 < cw_{ij} \le 1$ comprises

6   the concept weight, $0 < sw_{ij} \le 1$ comprises the structural weight, and $0 < rw_{ij} \le 1$

7   comprises the corpus weight for occurrence $j$ of concept $i$.

1        37.       (currently amended) A method according to Claim 36, further

2   comprising:

3            evaluating the concept weight in accordance with the formula:

4
$$cw_{ij} = \begin{cases} 0.25 + (0.25 \times t_{ij}), & 1 \le t_{ij} \le 3 \\ 0.25 + (0.25 \times [7 - t_{ij}]), & 4 \le t_{ij} \le 6 \\ 0.25, & t_{ij} \ge 7 \end{cases}$$

5   where $cw_{ij}$ comprises the concept weight and $t_{ij}$ comprises [[a]] the number of

6   terms for occurrence $j$ of each such concept $i$.

1        38.       (previously presented) A method according to Claim 36, further

2   comprising:

3            evaluating the structural weight in accordance with the formula:

4
$$sw_{ij} = \begin{cases} 1.0, & if(j \approx SUBJECT) \\ 0.8, & if(j \approx HEADING) \\ 0.7, & if(j \approx SUMMARY) \\ 0.5 & if(j \approx BODY) \\ 0.1 & if(j \approx SIGNATURE) \end{cases}$$

5   where $sw_{ij}$ comprises the structural weight for occurrence $j$ of each such concept $i$.

1      39.    (previously presented) A method according to Claim 36, further

2  comprising:

3      evaluating the corpus weight in accordance with the formula:

4
$$rw_{ij} = \begin{cases} \left(\dfrac{T - r_{ij}}{T}\right)^2, & r_{ij} > M \\ 1.0, & r_{ij} \leq M \end{cases}$$

5  where $rw_{ij}$ comprises the corpus weight, $r_{ij}$ comprises a reference count for

6  occurrence $j$ of each such concept $i$, $T$ comprises a total number of reference

7  counts of documents in the document set, and $M$ comprises a maximum reference

8  count of documents in the document set.

1      40.    (previously presented) A method according to Claim 36, further

2  comprising:

3      compressing the score in accordance with the formula:

4      $S_i' = \log(S_i + 1)$

5  where $S_i'$ comprises the compressed score for each such concept $i$.

1      41.    (original) A method according to Claim 35, further comprising:

2      maintaining concepts and terms in a global stop concept vector cache; and

3      filtering selection of the at least one concept based on the concepts and

4  terms maintained in the global stop concept vector cache.

1      42.    (original) A method according to Claim 35, further comprising:

2      identifying terms within at least one document in the document set; and

3      combining the identified terms into one or more of the concepts.

1      43.    (original) A method according to Claim 42, further comprising:

2      structuring each such identified term in the one or more concepts into

3  canonical concepts comprising at least one of word root, character case, and word

4  ordering.

1      44.    (original) A method according to Claim 42, further comprising:

2      including as terms at least one of nouns, proper nouns and adjectives.

1      Claims 45-47 (canceled).

1      48.    (previously presented) A method according to Claim 35, further

2      comprising:

3      calculating the similarity in accordance with the formula:

4
$$\cos \sigma_{AB} = \frac{\langle \vec{S}_A \cdot \vec{S}_B \rangle}{|\vec{S}_A||\vec{S}_B|}$$

5      where $\cos \sigma_{AB}$ comprises a similarity between a document $A$ and a document $B$,

6      $\vec{S}_A$ comprises a score vector for document $A$, and $\vec{S}_B$ comprises a score vector for

7      document $B$.

1      Claims 49-51 (canceled).

1      52.    (currently amended) A computer-readable storage medium holding

2      code for providing efficient document scoring of concepts within and clustering

3      of documents in an electronically-stored document set, comprising:

4      code for scoring a document in an electronically-stored document set,

5      comprising:

6      code for determining a frequency of occurrence of at least one

7      concept within a document;

8      code for analyzing a concept weight reflecting a specificity of

9      meaning for the at least one concept within the document, wherein the concept

10     weight is based on a number of terms for the at least one concept;

11     code for analyzing a structural weight reflecting a degree of

12     significance based on structural location within the document for the at least one

13     concept;

14          code for analyzing a corpus weight inversely weighing a reference

15    count of occurrences for the at least one concept within the document; and

16          code for evaluating a score to be associated with the at least one

17    concept as a function of a summation of the frequency, concept weight, structural

18    weight, and corpus weight;

19          code for forming the score assigned to the at least one concept as a

20    normalized score vector for each such document in the electronically-stored

21    document set;

22          code for determining a similarity between the normalized score vector for

23    each such document as an inner product of each normalized score vector;

24          code for grouping the documents by the score into a plurality of clusters,

25    comprising:

26          code for selecting a set of candidate seed documents from the

27    electronically-stored document set;

28          code for identifying seed documents by applying the similarity to

29    each such candidate seed document and selecting those candidate seed documents

30    that are sufficiently unique from other candidate seed documents as the seed

31    documents;

32          code for identifying a plurality of non-seed documents;

33          code for determining the similarity between each non-seed

34    document and a center of each cluster; and

35          code for assigning each non-seed document to the cluster with a

36    best fit, subject to a minimum fit; and

37          code for relocating outlier documents, comprising:

38          code for determining the similarity between each of the documents

39    grouped into each cluster based on the center of the cluster and the scores

40    assigned to each of the at least one concepts in that document;

41          code for dynamically determining a threshold for each cluster as a

42    function of the similarity between each of the documents; and

43      code for identifying and reassigning each of the documents with

44  the similarity falling outside the threshold.


1      53.      (currently amended) An apparatus for providing efficient

2  document scoring of concepts within and clustering of documents in an

3  electronically-stored document set, comprising:

4          means for scoring a document in an electronically-stored document set,

5  comprising:

6              means for determining a frequency of occurrence of at least one

7  concept within a document;

8              means for analyzing a concept weight reflecting a specificity of

9  meaning for the at least one concept within the document, wherein the concept

10  weight is based on a number of terms for the at least one concept;

11             means for analyzing a structural weight reflecting a degree of

12  significance based on structural location within the document for the at least one

13  concept;

14             means for analyzing a corpus weight inversely weighing a

15  reference count of occurrences for the at least one concept within the document;

16  and

17             means for evaluating a score to be associated with the at least one

18  concept as a function of a summation of the frequency, concept weight, structural

19  weight, and corpus weight;

20         means for forming the score assigned to the at least one concept as a

21  normalized score vector for each such document in the electronically-stored

22  document set;

23         means for determining a similarity between the normalized score vector

24  for each such document as an inner product of each normalized score vector;

25         means for grouping the documents by the score into a plurality of clusters,

26  comprising:

27             means for selecting a set of candidate seed documents from the

28  electronically-stored document set;

29           means for identifying seed documents by applying the similarity to

30    each such candidate seed document and selecting those candidate seed documents

31    that are sufficiently unique from other candidate seed documents as the seed

32    documents;

33           means for identifying a plurality of non-seed documents;

34           means for determining the similarity between each non-seed

35    document and a center of each cluster; and

36           means for assigning each non-seed document to the cluster with a

37    best fit, subject to a minimum fit; and

38          means for relocating outlier documents, comprising:

39           means for determining the similarity between each of the

40    documents grouped into each cluster based on the center of the cluster and the

41    scores assigned to each of the at least one concepts in that document;

42           means for dynamically determining a threshold for each cluster as

43    a function of the similarity between each of the documents; and

44           means for identifying and reassigning each of the documents with

45    the similarity falling outside the threshold.